# Data-Mining on GBytes of Encrypted Data

**Valeria Nikolaenko**

In collaboration with
Dan Boneh (Stanford); Udi Weinsberg, Stratis Ioannidis,
Marc Joye, Nina Taft (Technicolor).

# Outline

- **Motivation**
- Background on cryptographic tools
- Linear regression
- Our solution
- Experiments and performance
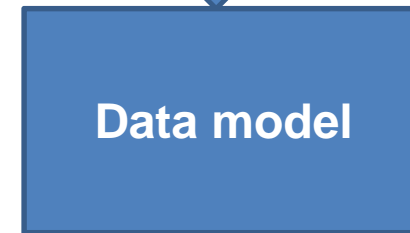
# Motivation

**Users**



Data

Data

**Data mining engine**

facebook   amazon.com   Google   NETFLIX

**Data model**

**Privacy concern!**

Engine learns **nothing** more than the model!

# Data Mining

- Classification
- Regression  : <span style="color:red">linear regression</span>
- Clustering
- Summarization  : <span style="color:red">matrix factorization</span>
- Dependency modeling

Main challenge: make these algorithms privacy preserving and efficient.

# Contribution

- Design of a practical system for privacy preserving linear regression
- Implementation
- Experiments on real datasets

Comparison to state of the art:
- Hall *et al.*'11: 2 days vs 3 min
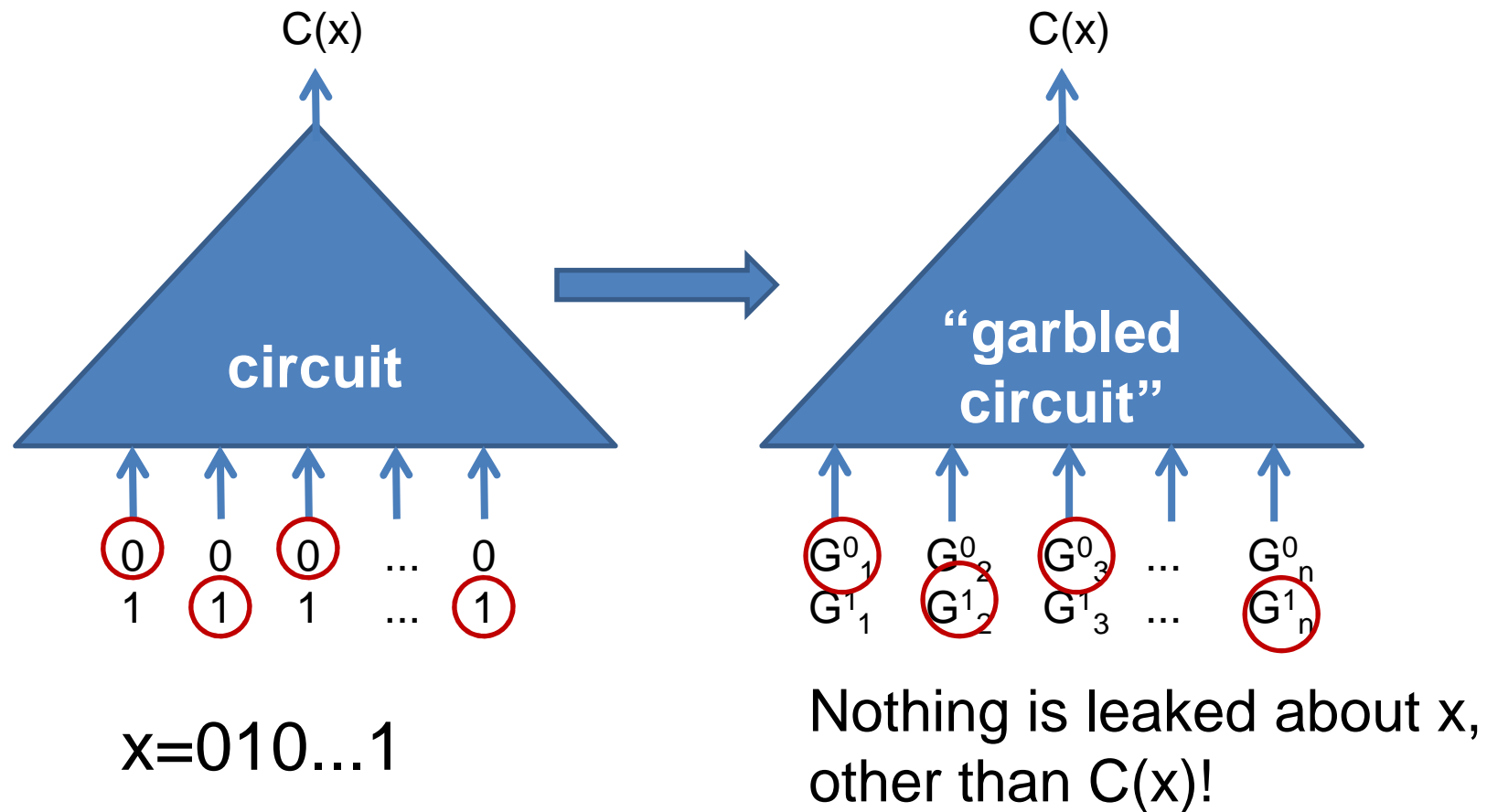- Graepel *et al.*'12: 10 min vs 2 sec

# Outline

- Motivation
- **Background on cryptographic tools**
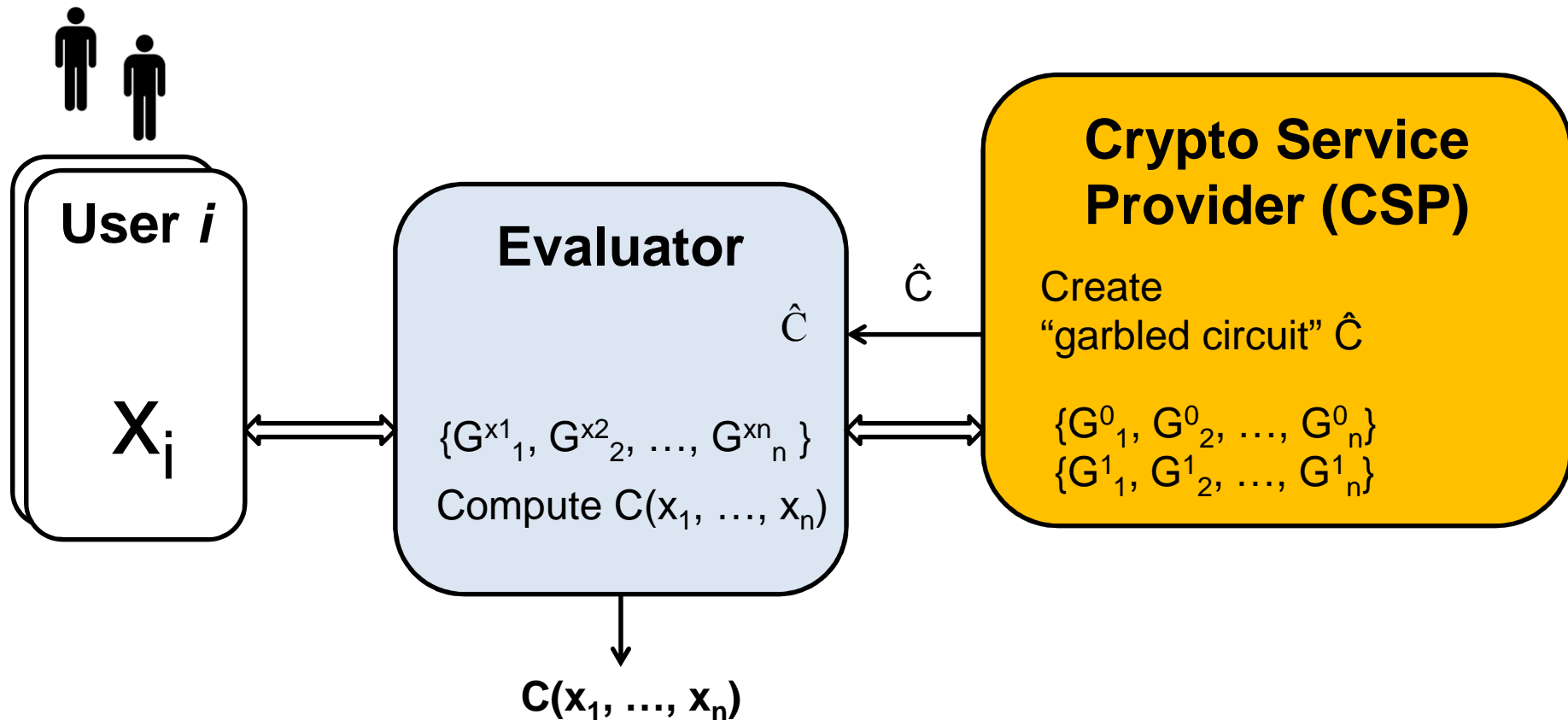- Linear regression
- Our solution
- Experiments and performance

# Computations on Encrypted Data

- 2009, **C. Gentry** – FHE
  (slow for our problems)
- 1979, **A. Shamir**
  1988, **BGW**          Secret sharing
  (huge communication overhead)
- 1982,  **A.C. Yao** – Garbled circuits

- Our approach: hybrid of Yao and hom. encryption

# Yao's Garbled Circuits



C(x)

circuit

⓪  0  ⓪  ...  0
1  ①  1  ...  ①

x=010...1

C(x)

"garbled circuit"

$G^0_1$  $G^0_2$  $G^0_3$  ...  $G^0_n$
$G^1_1$  $G^1_2$  $G^1_3$  ...  $G^1_n$

Nothing is leaked about x, other than C(x)!

# Data Mining System Architecture

**User *i***

$X_i$

**Evaluator**

$\hat{C}$

$\{G^{x1}_1, G^{x2}_2, \ldots, G^{xn}_n\}$

Compute $C(x_1, \ldots, x_n)$

$C(x_1, \ldots, x_n)$

**Crypto Service Provider (CSP)**

$\hat{C}$

Create "garbled circuit" $\hat{C}$

$\{G^0_1, G^0_2, \ldots, G^0_n\}$
$\{G^1_1, G^1_2, \ldots, G^1_n\}$

# System Properties

✓ Evaluator learns the model, not the inputs

Problems:

- Not scalable with the number of users
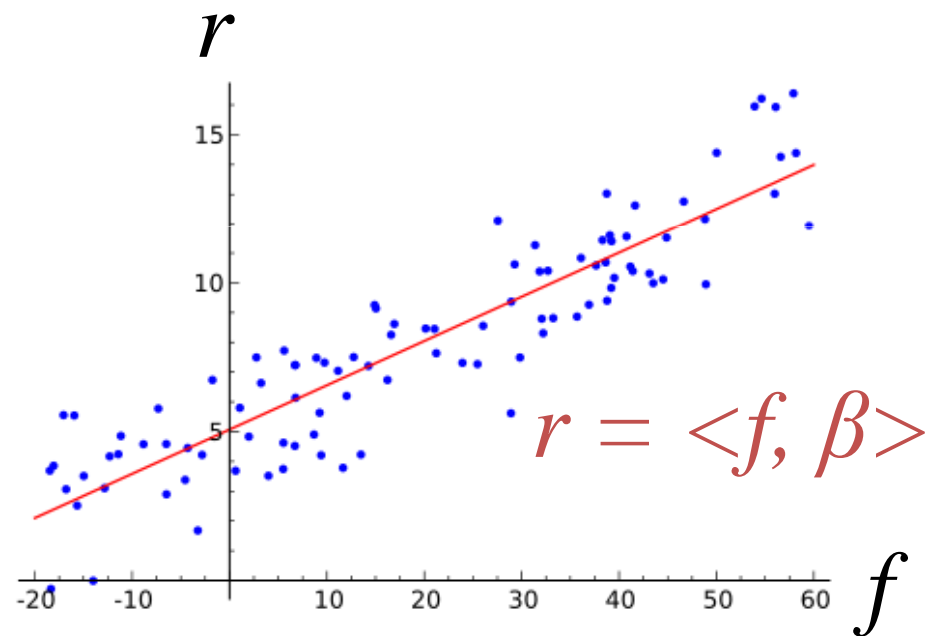- Users need to be online

# Outline

- Motivation
- Background on cryptographic tools
- **Linear regression**
- Our solution
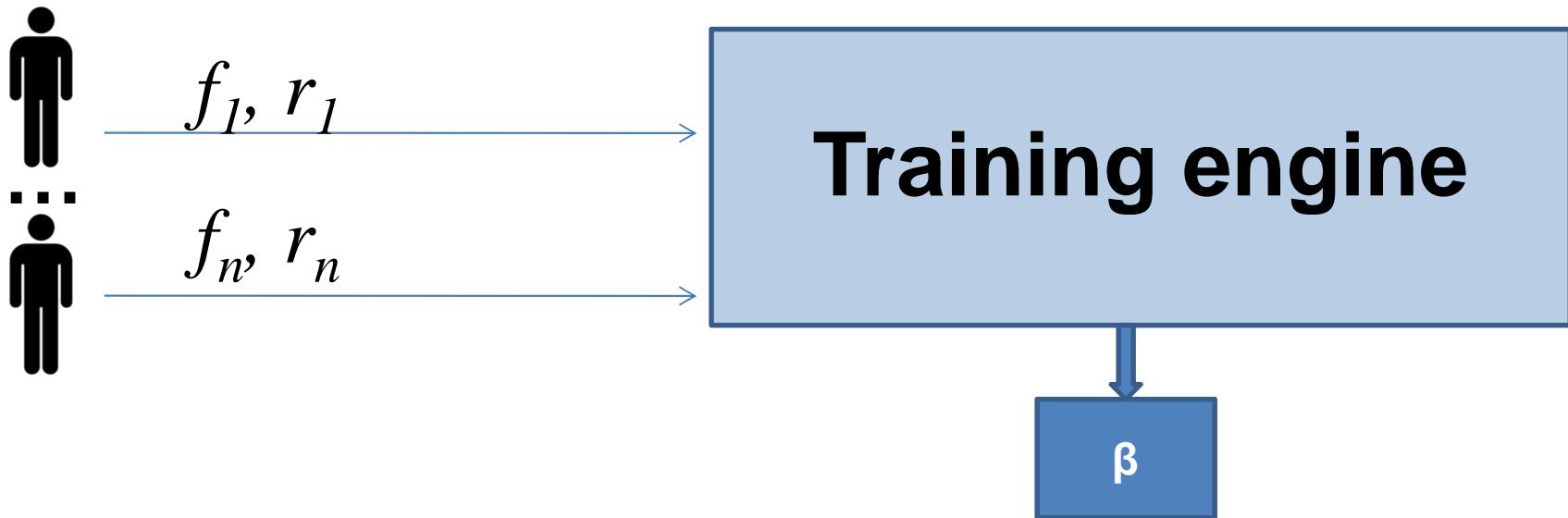- Experiments and performance

# Linear Regression I

- $(f, r)$ – from users
- solve for $\beta$

$$A\boldsymbol{\beta} = b$$



$r = <f, \beta>$

# Linear Regression II

**Users**



$f_1, r_1$

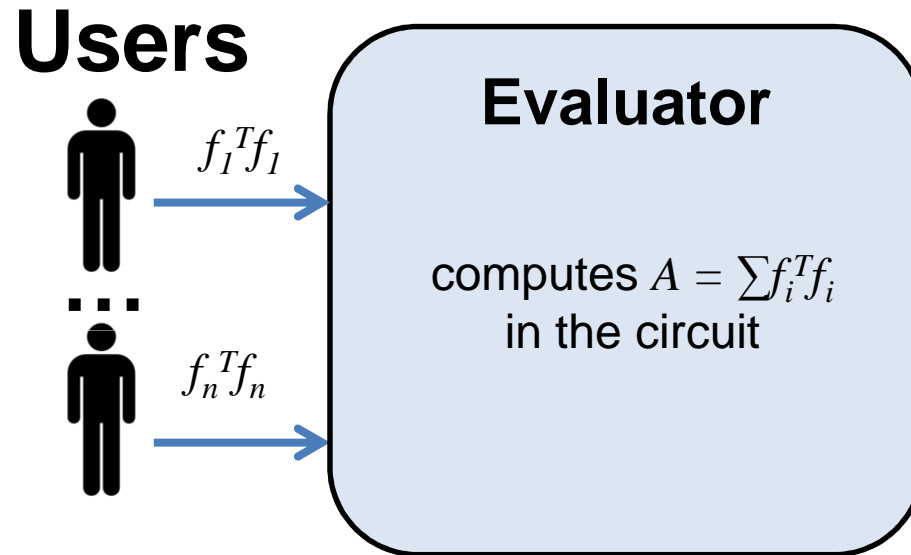$f_n, r_n$

**Training engine**

β

Training engine learns **nothing** about $f$'s and $r$'s, other than $\beta$!

# Outline

- Motivation
- Background on cryptographic tools
- Linear regression
- **Our solution**
- Experiments and performance

# Construct Matrices

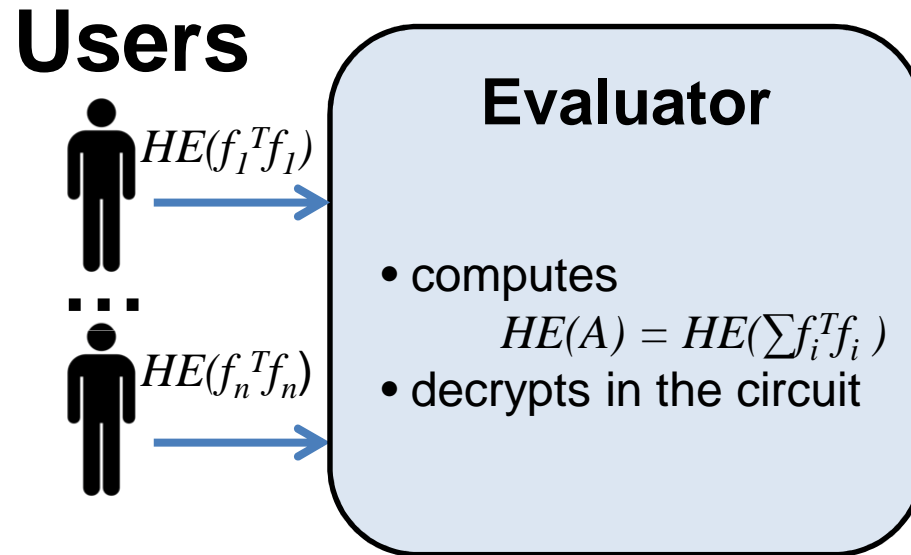Phase1: compute $A = \sum f_i^T f_i$ and $b = \sum f_i^T r_i$

**Users**



$f_1^T f_1$

**Evaluator**

computes $A = \sum f_i^T f_i$
in the circuit

$f_n^T f_n$

For additions can use <u>homomorphic encryption</u>:

$$[HE(A_1), HE(A_2)] \to HE(A_1 + A_2)$$

# Construct Matrices

Phase1: compute $A = \sum f_i^T f_i$ and $b = \sum f_i^T r_i$

**Users**



$HE(f_1^T f_1)$

...

$HE(f_n^T f_n)$

**Evaluator**

- computes
  $HE(A) = HE(\sum f_i^T f_i)$
- decrypts in the circuit

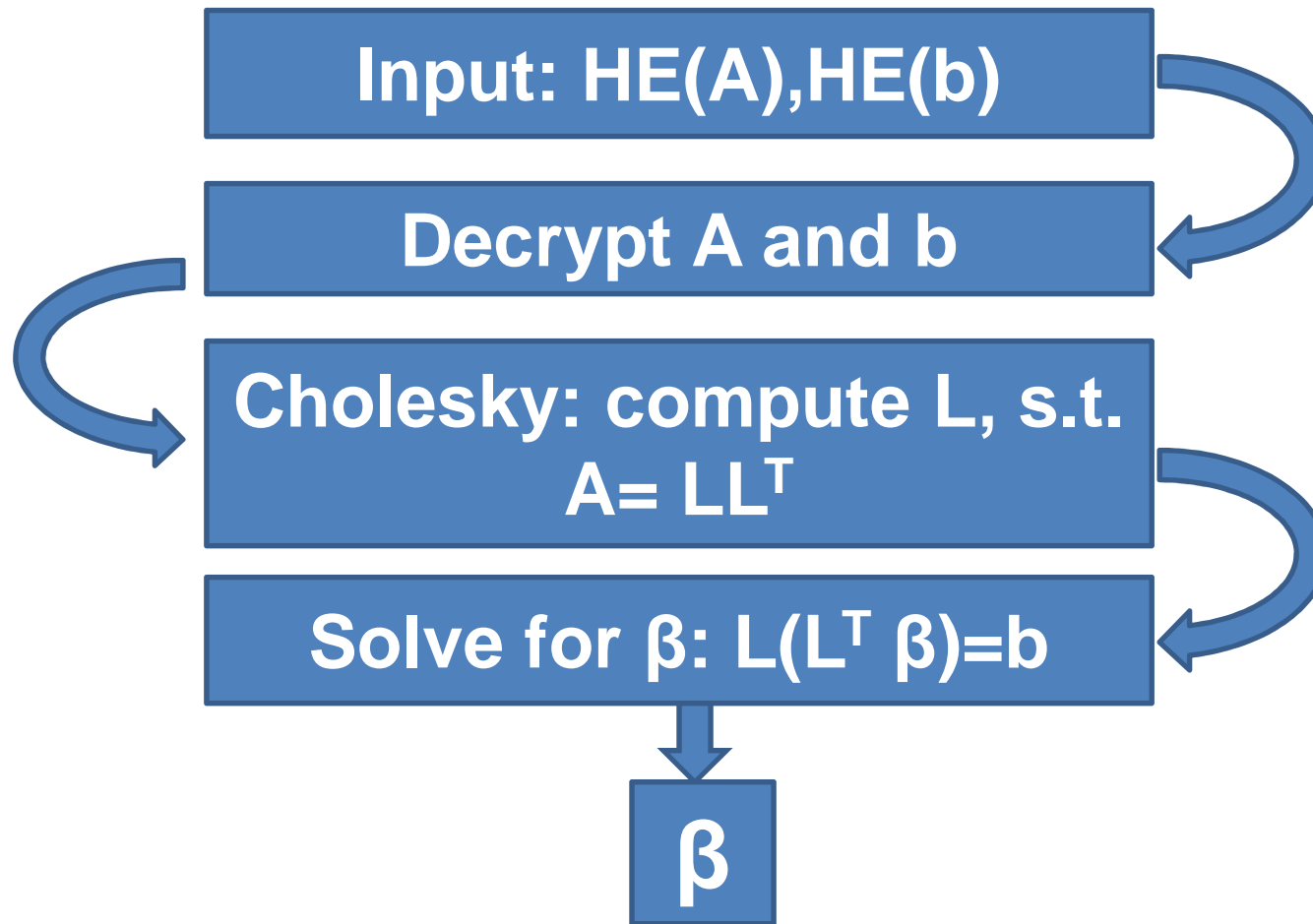For additions can use <u>homomorphic encryption</u>:

$$[HE(A_1), HE(A_2)] \rightarrow HE(A_1 + A_2)$$
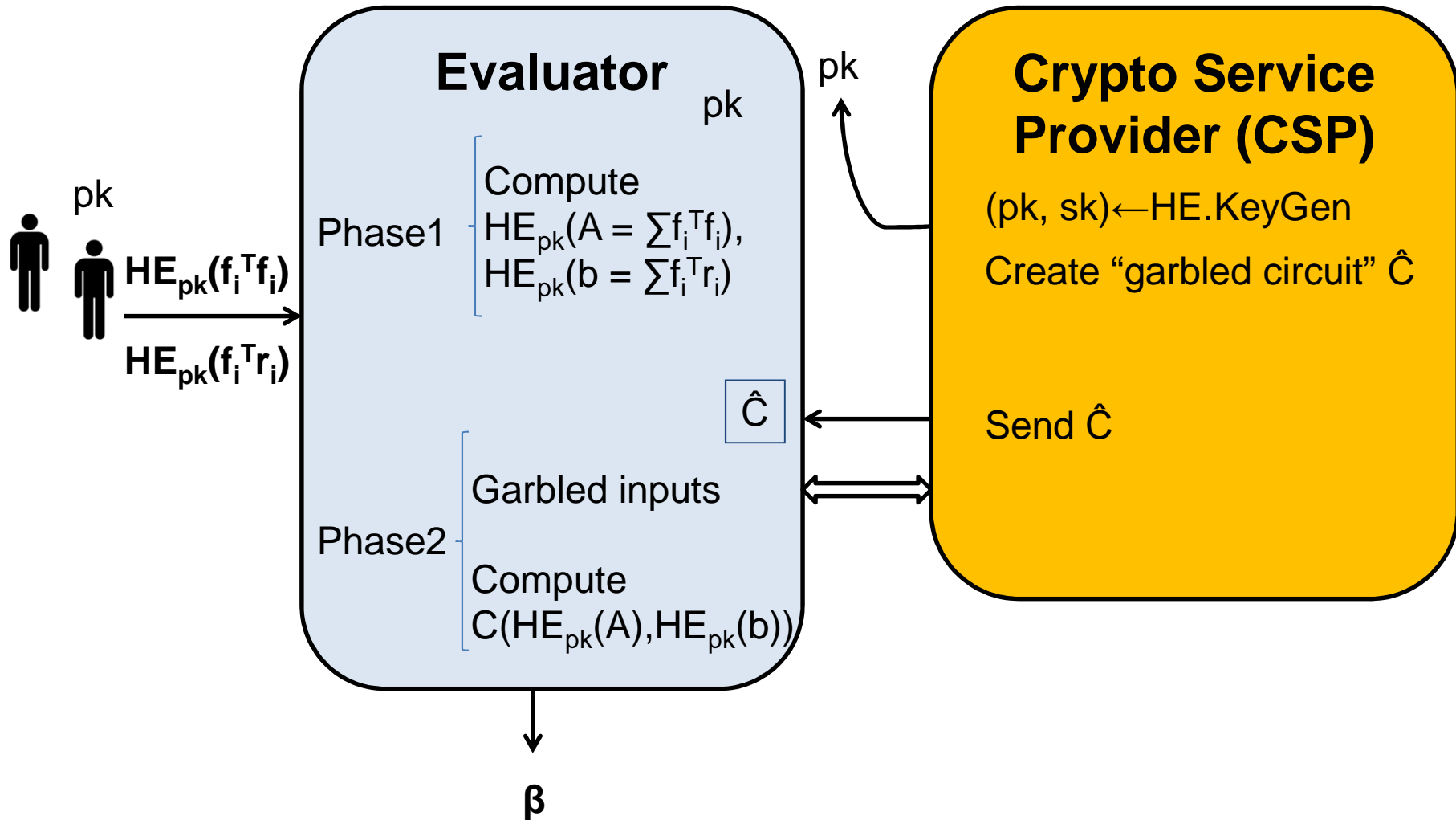
Circuit – independent of the number of users
Users – offline

# Solve Linear System

Phase2: solve for β  $A\beta = b$



Input: HE(A),HE(b)

Decrypt A and b

Cholesky: compute L, s.t. $A = LL^T$

Solve for β: $L(L^T \beta) = b$

β

# Privacy Preserving Regression System

# System Properties

✓ Evaluator learns the model, not the inputs

✓ Scalable with the number of users

✓ Users can be offline

## Extensions:

- Masking instead of decryption in circuit
- Protection against malicious

  Evaluator and CSP

# Outline

- Motivation
- Background on cryptographic tools
- Linear regression
- Our solution
- **Experiments and performance**

# Performance

- Hybrid vs. Pure-Yao: 100 times improvement in time!

- For up to 20 features, 1000's users
  - time < 3 min
  - communication < 1GB

- For 100 million of users, 20 features: 8.75 hours

- Tested on real datasets

# Conclusion

- Privacy-preserving data-mining is **efficient**
- Our approach can be used in practice

- Current work: matrix factorization

- Future work:
  - implement **other data mining algorithms**
  - improving implementation to support **high parallelization**

# Thank you!

Questions? valerini@stanford.edu